

<https://helda.helsinki.fi>

---

## Community-led, integrated, reproducible multi-omics with anvi'o

Eren, A. Murat

2021-01

---

Eren , A M , Kiefl , E , Shaiber , A , Veseli , I , Miller , S E , Schechter , M S , Fink , I , Pan , J N , Yousef , M , Fogarty , E C , Trigodet , F , Watson , A R , Esen , O C , Moore , R M , Clayssen , Q , Lee , M D , Kivenson , V , Graham , E D , Merrill , B D , Karkman , A , Blankenberg , D , Eppley , J M , Sjodin , A , Scott , J J , Vazquez-Campos , X , McKay , L J , McDaniel , E A , Stevens , S L R , Anderson , R E , Fuessel , J , Fernandez-Guerra , A , Maignien , L , Delmont , T O & Willis , A D 2021 , ' Community-led, integrated, reproducible multi-omics with anvi'o ' , Nature Microbiology , vol. 6 , pp. 3-6 . <https://doi.org/10.1038/s41564-020-00834-3>

---

<http://hdl.handle.net/10138/332078>

<https://doi.org/10.1038/s41564-020-00834-3>

---

acceptedVersion

---

*Downloaded from Helda, University of Helsinki institutional repository.*

*This is an electronic reprint of the original article.*

*This reprint may differ from the original in pagination and typographic detail.*

*Please cite the original version.*

# Community-led, integrated, reproducible multi-omics with anvio

A. Murat Eren<sup>1,2,3,\*</sup>, Evan Kiehl<sup>1,4</sup>, Alon Shaiber<sup>1,4</sup>, Iva Veseli<sup>1,4</sup>, Samuel E. Miller<sup>1</sup>, Matthew S. Schechter<sup>1,2</sup>, Isaac Fink<sup>1</sup>, Jessica N. Pan<sup>1</sup>, Mahmoud Yousef<sup>1</sup>, Emily C. Fogarty<sup>1</sup>, Florian Trigodet<sup>1</sup>, Andrea R. Watson<sup>1</sup>, Özcan C. Esen<sup>1</sup>, Ryan M. Moore<sup>5</sup>, Quentin Clayssen<sup>6</sup>, Michael D. Lee<sup>7,8</sup>, Veronika Kivenson<sup>9</sup>, Elaina D. Graham<sup>10</sup>, Bryan D. Merrill<sup>11</sup>, Antti Karkman<sup>12</sup>, Daniel Blankenberg<sup>13,14</sup>, John M. Eppley<sup>15</sup>, Andreas Sjödin<sup>16</sup>, Jarrod J. Scott<sup>17</sup>, Xabier Vázquez-Campos<sup>18</sup>, Luke J. McKay<sup>19,20</sup>, Elizabeth A. McDaniel<sup>21</sup>, Sarah L. R. Stevens<sup>22,23</sup>, Rika Anderson<sup>24</sup>, Jessika Fuessel<sup>1</sup>, Antonio Fernandez-Guerra<sup>25</sup>, Lois Maignien<sup>3,26</sup>, Tom O. Delmont<sup>27</sup>, Amy D. Willis<sup>28</sup>

<sup>1</sup> Department of Medicine, University of Chicago, Chicago, IL, USA; <sup>2</sup> Committee on Microbiology, University of Chicago, Chicago, IL, USA; <sup>3</sup> Josephine Bay Paul Center for Comparative Molecular Biology and Evolution, Marine Biological Laboratory, Woods Hole, MA, USA; <sup>4</sup> Graduate Program in Biophysical Sciences, University of Chicago, Chicago, IL, USA; <sup>5</sup> Center for Bioinformatics and Computational Biology, University of Delaware, DE, USA; <sup>6</sup> Department of Biology, Institute of Microbiology, ETH Zurich, Zurich, Switzerland; <sup>7</sup> Exobiology Branch, NASA Ames Research Center, Mountain View, CA, USA; <sup>8</sup> Blue Marble Space Institute of Science, Seattle, WA, USA; <sup>9</sup> Department of Microbiology, Oregon State University, Corvallis, OR, USA; <sup>10</sup> Department of Biological Sciences, University of Southern California, CA, USA; <sup>11</sup> Department of Microbiology and Immunology, Stanford University School of Medicine, Stanford, CA, USA; <sup>12</sup> Department of Microbiology, University of Helsinki, Helsinki, Finland; <sup>13</sup> Genomic Medicine Institute, Lerner Research Institute, Cleveland Clinic, Cleveland, OH, USA; <sup>14</sup> Department of Molecular Medicine, Cleveland Clinic Lerner College of Medicine, Case Western Reserve University, Cleveland, OH, USA; <sup>15</sup> Daniel K. Inouye Center for Microbial Oceanography: Research and Education, University of Hawaii, Manoa, Honolulu, HI, USA; <sup>16</sup> Division of CBRN Security and Defence, Swedish Defence Research Agency - FOI, Umeå, Sweden; <sup>17</sup> Smithsonian Tropical Research Institute, Bocas del Toro, Republic of Panamá; <sup>18</sup> NSW Systems Biology Initiative, School of Biotechnology and Biomolecular Sciences, The University of New South Wales, Sydney, NSW 2052, Australia; <sup>19</sup> Center for Biofilm Engineering, Montana State University, Bozeman, MT, USA; <sup>20</sup> Department of Land Resources and Environmental Sciences, Montana State University, Bozeman, MT, USA; <sup>21</sup> Department of Bacteriology, University of Wisconsin-Madison, Madison, WI, USA; <sup>22</sup> Wisconsin Institute for Discovery, University of Wisconsin-Madison, Madison, WI, USA; <sup>23</sup> American Family Insurance Data Science Institute, University of Wisconsin-Madison, Madison, WI, USA; <sup>24</sup> Department of Biology, Carleton College, Northfield, MN, USA; <sup>25</sup> Lundbeck GeoGenetics Centre, The Globe Institute, University of Copenhagen, 1350 Copenhagen, Denmark; <sup>26</sup> Laboratoire de Microbiologie des Environnements Extrêmes (LM2E), Univ Brest, CNRS, Ifremer, Plouzané, France; <sup>27</sup> Génomique Métabolique, Genoscope, Institut François Jacob, CEA, CNRS, Univ Evry, Université Paris-Saclay, 91057 Evry, France; <sup>28</sup> Department of Biostatistics, University of Washington, Seattle, WA, USA.

\* Correspondence: [meren@uchicago.edu](mailto:meren@uchicago.edu)

## Standfirst

**Big data abounds in microbiology, but the workflows designed to enable researchers to interpret data can constrain the biological questions that can be asked. Five years after *anvi'o* was first published, this community-led multi-omics platform is maturing into an open software ecosystem that reduces constraints in 'omics data analyses.**

Generating hundreds of millions of sequences from a microbial habitat is now commonplace for many microbiologists<sup>1</sup>. While the massive data streams offer detailed snapshots of the lifestyles of microorganisms, this data revolution in microbiology means that a new generation of computational tools is needed to empower life scientists in the era of multi-omics.

To meet the growing computational needs of the life sciences, computer scientists and bioinformaticians have created thousands of software tools<sup>2</sup>. These software fall into two general categories: 'essential tools' that implement functions fundamental to most bioinformatics tasks, and 'workflows' that make specific analytic strategies accessible.

If a comprehensive microbial 'omics investigation is a sophisticated dish, then essential tools are the kitchenware needed to cook. A chef can combine them in unique ways to answer any question, yet such freedom in data analysis not only requires the mastery of each essential tool but also demands experience in data wrangling and fluency in the command line environment to match the output format of one tool to the input requirements of the next. This barrier is overcome by workflows, which implement popular analysis strategies and make them accessible to those who have limited training in computation. If a comprehensive microbial 'omics investigation is a sophisticated dish, then each 'omics workflow is a recipe that turns raw material into a specific meal. For instance, a workflow for 'pangenomics' would typically take in a set of genomes and (1) identify open reading frames in all input genomes, (2) reciprocally align all translated amino acid sequences, (3) identify gene clusters by resolving pairwise sequence homology across all genes, and (4) report the distribution of gene clusters across genomes. By doing so, a software that implements pangenomics, such as Roary<sup>3</sup>, would seamlessly run multiple essential tools consecutively, resolve input/output requirements of each, and address various ad hoc computational challenges to concoct a pangenome. Popular

efforts to make accessible workflows that form the backbone of 'omics-based microbiological studies include the Galaxy platform<sup>4</sup>, bioBakery software collection<sup>5</sup>, M-Tools (i.e., GroopM<sup>6</sup>, CheckM<sup>7</sup>), and KBase<sup>8</sup>. While 'omics workflows conveniently summarise raw data into tables and figures, the ability to analyse data beyond pre-defined strategies they implement continues to be largely limited to master chefs, presenting the developers of 'omics workflows with a substantial responsibility: pre-determining the investigative routes their software enables users to traverse, which can influence how researchers interact with their data, conceivably affecting biological interpretations.

We introduced anvi'o (an analysis and visualisation platform for 'omic data) as an alternative solution for microbiologists who wanted more freedom in research questions they could ask of their data<sup>9</sup>. We started with what we regarded as the most pressing need at the time: a platform that enabled the reconstruction and interactive refinement of microbial genomes from environmental metagenomes. Fundamentals of this strategy were already established by those who pioneered genome-resolved metagenomics<sup>10</sup>, but interactive visualisation and editing software that would enable microbiologists to intimately work with metagenome-assembled genomes was lacking. During the past five years anvi'o has become a community-driven software platform that currently stands upon more than 90,000 lines of open-source code and supports interactive and fully integrated access to state-of-the-art 'omics strategies including genomics, genome-resolved metagenomics and metatranscriptomics, pangenomics, metapangenomics, phylogenomics, and microbial population genetics (Figure 1).

Anvi'o differs from existing bioinformatics software due to its modular architecture, which enables flexibility, interactivity, reproducibility, and extensibility. To achieve this, the platform contains more than 100 interoperable programs, each of which performs individual tasks that can be combined to build new and unique analytical workflows. Anvi'o programs generate, modify, query, split, and merge anvi'o projects, which are really a set of extensible, self-contained SQLite databases. The interconnected nature of anvi'o programs which are glued together by these common data structures yields a network (<http://merenlab.org/nt>), rather than predetermined, linear paths for analysis. Through this modularity, anvi'o empowers its users to navigate through 'omics data without imposing rigid workflows.

Integrated interactive visualisation is at the center of anvi'o and helps researchers to engage with their data in all stages of analysis. Within the same interface, an anvi'o user can visualise amino acid sequence alignments between homologous genes across multiple genomes, investigate nucleotide-level coverage patterns and variants on the same DNA segment across metagenomes, interrogate associations between the genomic abundance and transcriptomic activity of environmental microbes, display phylogenetic trees and clustering dendrograms, and more. Furthermore, users can extend anvi'o displays with project-specific external data, increasing the utility of interactive interfaces for holistic descriptions of complex systems. The anvi'o interactive interface also provides its users with the artistic freedom to change colours, sizes, and drawing styles of display objects, add annotations, or reorder data layers for detailed communication of intricate observations. Because each anvi'o project is self-contained, researchers can easily make their analyses available online as a whole or in part, thereby enabling the integration, reusability, and reproducibility of their findings beyond static figures or tables. This strategy promotes transparency by permitting community validation and scrutiny through full access to data that underlie final conclusions.

Several key studies that used anvi'o during the past few years have demonstrated the integrative capabilities of the platform by implementing a combination of 'omics strategies to facilitate in-depth analysis of naturally occurring microbial habitats. For instance, Reveillaud and Bordenstein et al. reconstructed new genomes of *Wolbachia*, a fastidious endosymbiont<sup>11</sup>, from individual insect ovary metagenomes, and computed a pangenome to compare these novel genomes to an existing reference<sup>12</sup>. They were then able to characterise the ecology of gene clusters in the environment by effectively combining metagenomics and pangenomics, discovering new members of the *Wolbachia* mobilome<sup>12</sup>. Yeoman et al. combined phylogenomics and pangenomics to infer ancestral relationships between a set of cultivar and metagenome-assembled genomes through a *de novo* identified set of single-copy core genes<sup>13</sup>. They demonstrated the correspondence among these genomes based on gene cluster membership patterns, phylogenomic inference, and average nucleotide identity in a single display<sup>13</sup>. Delmont and Kiefl et al. characterised the population structure of a subclade of SAR11, one of the most abundant microbial populations on Earth, by describing the

environmental core genes of a single genome across surface ocean metagenomes<sup>14</sup>. By linking single-amino acid variants in the environment to the predicted tertiary structures of these genes, they combined microbial population genetics with protein biochemistry to shed light on distinct evolutionary processes shaping the population structures of these bacteria<sup>14</sup>. Each of these studies employs unique approaches beyond well-established 'omics workflows to create rich, reproducible, and shareable data products (see <http://merenlab.org/data>).

Anvi'o does not implement strategies that take in raw data and produce summary tables or figures via a single command. As a result, anvi'o has a relatively steep learning curve. To address this, we have written extensive online tutorials that currently exceed 120,000 words, organised free workshops for hands-on anvi'o training, and created open educational resources to learn microbial 'omics. To interact with anvi'o users we set up an online forum and messaging service. During the past two years, more than 750 registered members of these services have engaged in technical and scientific discussions via more than 9,000 messages. But even when resources for learning are available, the journey from raw 'omics data to biological insights often takes a significant number of atomic steps of computation. To ameliorate the burden of scale and reproducibility in big data analyses we have also introduced anvi'o workflows, which automate routine computational steps of commonly used analytical strategies in microbial 'omics (<http://merenlab.org/anvio-workflows>). The anvi'o workflows are powered by Snakemake<sup>15</sup>, which ensures relatively easy deployment to any computer system and automatic parallelisation of independent analysis steps. By turning raw input into data products to be analysed in the anvi'o software ecosystem, anvi'o workflows reduce the barriers for advanced use of computational resources and processing of large data streams for microbial 'omics.

As the developers of anvi'o who strive to create an open community resource, our next big challenge is to attract bioinformaticians to consider anvi'o as a software development ecosystem they can use for their own science. Any program that reads from or writes to anvi'o projects either directly (in any modern programming language) or through anvi'o application programmer interfaces (in Python) will immediately become accessible to anvi'o users, and

such applications will benefit from the data integration, interactive data visualisation, and error checking assurances anvi'o offers.

As an open-source platform that empowers microbiologists by offering them integrated yet uncharted means to steer through complex 'omics data, anvi'o welcomes its new users and contributors.

## Acknowledgements

The URL <https://github.com/merenlab/anvio/blob/master/AUTHORS.txt> serves a complete list of anvi'o developers. We thank the creators of other open-source software tools for their generosity, anvi'o users for their patience with us, and Karen Lolans (0000-0003-1903-756X) for her critical reading of the manuscript and suggestions. The authors gratefully acknowledge support for anvi'o from the Simons Foundation and Alfred P. Sloan Foundation.

## Author contributions

AME, EK, AS, IV, SEM, MSS, IF, JNP, MY, ECF, FT, ARW, OCE, RMM, QC, and ADW coded and documented anvi'o, contributed to the implementation of new analytical strategies, and engaged with the anvi'o community. MDL, VK, EDG, BDM, and AK wrote blog posts and tutorials to make anvi'o accessible to the broader community. XVC, LJM helped with technical issues and testing of new features on GitHub. EAM, SLRS, and RA created undergraduate and graduate-level educational material and taught anvi'o. LM organized workshops for the training of research professionals. JF, AFG, LM, TOD, and ADW made intellectual contributions that influenced the direction of the platform. AME wrote the paper and prepared the figure with input from all authors.

## Competing interests

Authors have no conflicts of interest to declare.

## Figure Legends

**Figure 1.** A glimpse of the interconnected nature of 'omics analysis strategies anvi'o makes accessible, and their potential applications.

## References

1. White, R. A., Callister, S. J., Moore, R. J., Baker, E. S. & Jansson, J. K. The past, present and future of microbiome analyses. *Nat. Protoc.* **11**, 2049–2053 (2016).
2. Callahan, A., Winnenburg, R. & Shah, N. H. U-Index, a dataset and an impact metric for informatics tools and databases. *Sci Data* **5**, 180043 (2018).
3. Page, A. J. *et al.* Roary: rapid large-scale prokaryote pan genome analysis. *Bioinformatics* **31**, 3691–3693 (2015).
4. Jalili, V. *et al.* The Galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2020 update. *Nucleic Acids Res.* **48**, W395–W402 (2020).
5. McIver, L. J. *et al.* bioBakery: a meta'omic analysis environment. *Bioinformatics* vol. 34 1235–1237 (2018).
6. Imelfort, M. *et al.* GroopM: an automated tool for the recovery of population genomes from related metagenomes. *PeerJ* **2**, e603 (2014).
7. Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM: assessing the quality of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome Res.* **25**, 1043–1055 (2015).
8. Arkin, A. P. *et al.* KBase: The United States Department of Energy Systems Biology Knowledgebase. *Nat. Biotechnol.* **36**, 566–569 (2018).
9. Eren, A. M. *et al.* Anvi'o: an advanced analysis and visualization platform for 'omics data. *PeerJ* **3**,



198 e1319 (2015).

199 10. Tyson, G. W. *et al.* Community structure and metabolism through reconstruction of microbial  
200 genomes from the environment. *Nature* **428**, 37–43 (2004).

201 11. Werren, J. H., Baldo, L. & Clark, M. E. Wolbachia: master manipulators of invertebrate biology. *Nat.*  
202 *Rev. Microbiol.* **6**, 741–751 (2008).

203 12. Reveillaud, J. *et al.* The Wolbachia mobilome in *Culex pipiens* includes a putative plasmid. *Nat.*  
204 *Commun.* **10**, 1051 (2019).

205 13. Yeoman, C. J. *et al.* Genome-resolved insights into a novel Spiroplasma symbiont of the Wheat  
206 Stem Sawfly (*Cephus cinctus*). *PeerJ* **7**, e7548 (2019).

207 14. Delmont, T. O. *et al.* Single-amino acid variants reveal evolutionary processes that shape the  
208 biogeography of a global SAR11 subclade. *Elife* **8**, (2019).

209 15. Köster, J. & Rahmann, S. Snakemake--a scalable bioinformatics workflow engine. *Bioinformatics* **28**,  
210 2520–2522 (2012).

